

Machine Learning

Noel Moreno Lemus, Ph.D.

Table of contents

Introduction	1
Overview of Machine Learning:	2
Applications of Machine Learning	2
Key Concepts and Terminology	3
Versions Used in this Book	4
 I. Fundamentals	 5
1. Fundamentals of Machine Learning	7
1.1. Probability and Statistics	7
1.1.1. Introduction to Probability	7
1.2. Linear Algebra	9
1.3. Optimization	9
1.4. Data Preprocessing and Feature Engineering	9
 II. Optimization	 11
2. Optimization	13
2.1. Integer Programming	20
2.2. Other topics	20

Table of contents

III. Supervised Learning	21
3. Linear Regression	25
3.1. Learning by Examples	26
4. Logistic Regression	29
5. Decision Trees and Random Forests	31
6. Support Vector Machines	33
7. Neural Networks	35
IV. Unsupervised Learning	37
8. Unsupervised Learning	39
8.1. Clustering	39
8.2. Dimensionality Reduction	39
8.3. Anomaly Detection	39
V. Deep Learning	41
9. Neural Networks and Backpropagation	43
10. Convolutional Neural Networks (CNNs)	45
11. Recurrent Neural Networks (RNNs)	47
12. Applications in Computer Vision and Natural Language Processing	49

Table of contents

VI. Advanced Topics	51
13. Graph Neural Networks	53
14. Reinforcement Learning	55
15. Generative Models	57
16. Model Interpretability	59
17. Conclusion	61
References	63
Appendices	63
A. Appendices	65
A.1. Mathematical Proofs	65
A.2. Additional Resources	65

Introduction

Work in Progress

This is a Work in Progress project. The idea is to summarize all ML topics and how I see and understand them.

Machine learning is a rapidly growing field that enables computers to learn from data, without being explicitly programmed. The goal of machine learning is to build models that can make predictions or take actions based on input data, and improve their performance over time through experience.

Note

Note that there are five types of callouts, including: **note**, **warning**, **important**, **tip**, and **caution**.

Tip With Caption

This is an example of a callout with a caption.

Expand To Learn About Collapse

This is an example of a ‘folded’ caution callout that can be expanded by the user. You can use `collapse="true"` to collapse it by default or `collapse="false"` to make a collapsible callout that is expanded

by default.

Overview of Machine Learning:

Machine learning is a subfield of artificial intelligence that involves the development of algorithms and statistical models that allow computers to learn from data. There are three main types of machine learning: supervised learning, unsupervised learning, and reinforcement learning.

Supervised learning is the most common type of machine learning, in which a model is trained on a labeled dataset to make predictions about new, unseen data. Examples include linear regression, logistic regression, and decision trees.

Unsupervised learning involves discovering patterns in unlabeled data, such as clustering and dimensionality reduction.

Reinforcement learning involves training an agent to make decisions in an environment to maximize a reward.

Applications of Machine Learning

Machine learning has many applications in various industries, including:

- Healthcare: for example, identifying potential health risks, diagnosing diseases, and creating personalized treatment plans
- Finance: for example, detecting fraudulent transactions, predicting stock prices, and identifying potential investment opportunities
- Retail: for example, personalizing product recommendations, optimizing pricing strategies, and improving supply chain efficiency
- Manufacturing: for example, predictive maintenance, quality control, and optimization of production processes

Key Concepts and Terminology

- Transportation: for example, traffic prediction, autonomous driving, and fleet management
- Cybersecurity: for example, intrusion detection, anomaly detection, and threat intelligence

Key Concepts and Terminology

Machine learning is a complex field with many technical terms and concepts. Some key terms and concepts that will be covered in this book include:

- Model: a representation of the relationships between input data and output predictions or actions
- Training: the process of fitting a model to a dataset
- Testing: the process of evaluating a model on new, unseen data
- Overfitting: when a model is too complex and performs well on the training data but poorly on the test data
- Regularization: a technique for preventing overfitting by adding a penalty term to the model's objective function
- Gradient descent: an optimization algorithm for finding the minimum of a function
- Neural networks: a type of model that is inspired by the structure and function of the human brain
- Convolutional neural networks (CNNs): a type of neural network designed for image recognition
- Recurrent neural networks (RNNs): a type of neural network designed for sequential data such as time series and natural language.

Versions Used in this Book

```
import sys
print("Python version: {}".format(sys.version))
import pandas as pd
print("pandas version: {}".format(pd.__version__))
import matplotlib
print("matplotlib version: {}".format(matplotlib.__version__))
import numpy as np
print("NumPy version: {}".format(np.__version__))
import scipy as sp
print("SciPy version: {}".format(sp.__version__))
import IPython
print("IPython version: {}".format(IPython.__version__))
import sklearn
print("scikit-learn version: {}".format(sklearn.__version__))
```

```
Python version: 3.9.2 (default, Feb 28 2021, 17:03:44)
[GCC 10.2.1 20210110]
pandas version: 1.5.2
matplotlib version: 3.6.3
NumPy version: 1.24.1
SciPy version: 1.10.0
IPython version: 8.8.0
```

```
scikit-learn version: 1.2.0
```

Part I.

Fundamentals

1. Fundamentals of Machine Learning

1.1. Probability and Statistics

1.1.1. Introduction to Probability

1.1.1.1. Definition of probability

Probability is a measure of the likelihood of an event occurring. It is a value between 0 and 1, where 0 indicates that an event will never occur and 1 indicates that an event will always occur.

Probability can be defined in different ways, but one of the most common ways is through the use of relative frequency. If we repeat an experiment many times and count the number of times an event of interest occurs, we can calculate the probability of that event as the ratio of the number of successful outcomes to the total number of trials. For example, if we flip a coin 10 times and it comes up heads 6 times, we can say that the probability of getting heads is $6/10$ or 0.6.

Probability can also be defined through the use of theoretical models. For example, in the coin flipping example, we can assume that the coin is fair and that the probability of getting heads is 0.5.

Probability can be applied to many different types of events and situations, such as in gambling, finance, weather forecasting, medical diagnosis,

1. *Fundamentals of Machine Learning*

and many more. In machine learning, probability is used to model the uncertainty of predictions, estimate model parameters and evaluate model performance.

1.1.1.2. **Random variables and events**

A random variable is a variable that takes on different values based on the outcome of a random experiment. The values of a random variable can be numerical or categorical, and the probability of each value is defined by a probability distribution.

For example, in a coin-tossing experiment, the random variable X can take on the values of “heads” or “tails”. The probability of getting heads is 0.5, and the probability of getting tails is also 0.5. We can represent the probability distribution of X in a table or a graph.

An event is a set of outcomes from a random experiment. For example, in a coin-tossing experiment, the event “getting heads” is the set $\{\text{heads}\}$, and the event “getting tails” is the set $\{\text{tails}\}$.

A random variable is said to be discrete if it can take on only a countable number of values and continuous if it can take on any value in an interval.

For example, in a dice-rolling experiment, the random variable X can take on the values 1, 2, 3, 4, 5, or 6. X is a discrete random variable.

On the other hand, in a temperature measurement experiment, the random variable X can take on any value between -273.15 and infinity (the absolute zero and the maximum temperature). X is a continuous random variable.

In machine learning, random variables are used to represent the input and output of a model, the parameters of a model and the noise in the data. Understanding the properties of random variables and the events they can generate is important to design and analyze machine learning algorithms.

1.2. Linear Algebra

1.1.1.3. Sample space and event space

1.1.1.4. Axioms of probability

1.2. Linear Algebra

1.3. Optimization

1.4. Data Preprocessing and Feature Engineering

Part II.

Optimization

2. Optimization

Optimization can mean different things. Mathematically speaking, it is possible to write an optimization problem in the form of:

$$\begin{aligned} & \min_{\mathbf{x} \in \mathbb{R}^n} f_i(\mathbf{x}), \quad (i = 1, 2, \dots, M) \\ & \text{subject to } \phi_j(\mathbf{x}) = 0, \quad (j = 1, 2, \dots, J) \\ & \psi_k(\mathbf{x}) \leq 0, \quad (k = 1, 2, \dots, K) \end{aligned}$$

$f_i(\mathbf{x})$, $\phi_j(\mathbf{x})$ and $\psi_k(\mathbf{x})$ are functions of the vector

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

where the components x_i of \mathbf{x} are called decision variables and they can be real continuous, discrete or a mixture of both.

The functions $f_i(\mathbf{x})$, ($i = 1, 2, \dots, M$) are called the objective functions, and in the case of $M = 1$ there is only one objective. The space spanned by the decision variables is called the search space \mathbb{R}^n , while the space formed by the objective function values is called the solution space. The objective function can be linear or nonlinear.

The equalities $\phi_j(\mathbf{x})$ and inequalities $\psi_k(\mathbf{x})$ are called constraints. The simplest case of the constraints is where x_i is $x_{i,\min} \leq x_i \leq x_{i,\max}$, which is called bounds.

2. Optimization

If the functions $f_i(\mathbf{x})$, $\phi_j(\mathbf{x})$ and $\psi_k(\mathbf{x})$ are all linear, then we are in presence of a linear programming problem. For linear programming problems, a significant progress was the development of the simplex method in 1947 by George B. Dantzig.

A special class of optimization is where there is no constraint at all and the only task is to find the minimum or maximum of a single objective function $f_i(\mathbf{x})$.

For example, imagine that you want to find the minimum of the Rosenbrock banana function:

$$f(x, y) = (1 - x)^2 + 100 \times (y - x^2)^2$$

```
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import cm
from matplotlib.ticker import LinearLocator, FormatStrFormatter
from mpl_toolkits.mplot3d import Axes3D

b = 10
f = lambda x,y: (1-x)**2 + b*(y-x**2)**2

# Initialize figure
figRos = plt.figure(figsize=(12, 7))
axRos = figRos.add_subplot(projection='3d')

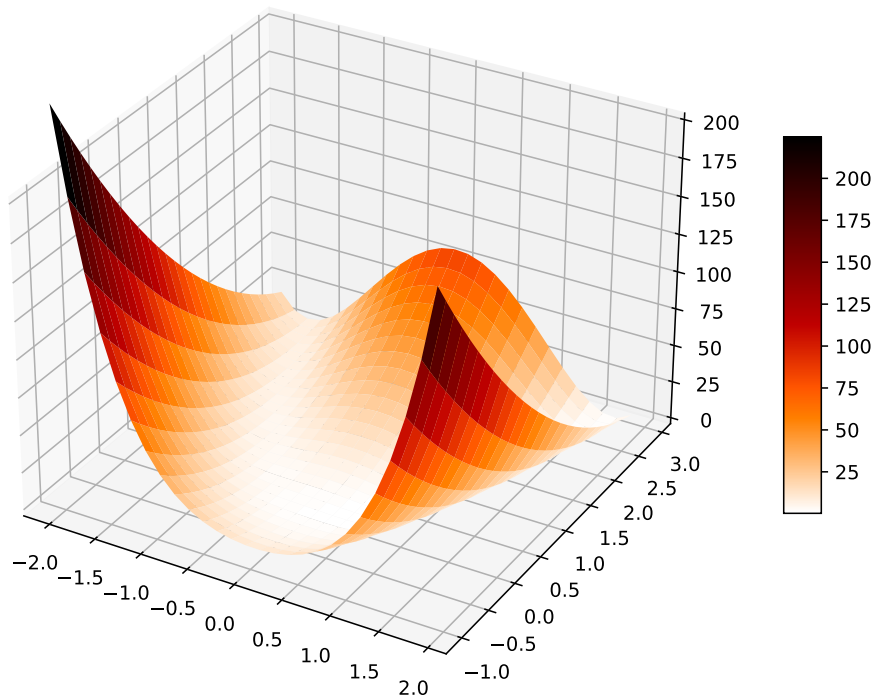
# Evaluate function
X = np.arange(-2, 2, 0.15)
Y = np.arange(-1, 3, 0.15)
X, Y = np.meshgrid(X, Y)
Z = f(X,Y)

# Plot the surface
```

```

surf = axRos.plot_surface(X, Y, Z, cmap=cm.gist_heat_r,
                          linewidth=0, antialiased=False)
axRos.set_zlim(0, 200)
figRos.colorbar(surf, shrink=0.5, aspect=10)
plt.show()

```



In order to find its minimum, we can set its partial derivatives to zero:

$$\frac{\partial f}{\partial x} = 2(1 - x) - 400(y - x^2)x = 0$$

2. Optimization

$$\frac{\partial f}{\partial y} = 200(y - x^2) = 0$$

The second equation implies that $y = x^2$ and substituted on the first one we get $x = 1$ as a solution. Then the f_{min} occurs at $x = y = 1$.

```
plt.figure(figsize=(8, 5))
plt.contour(X,Y,Z,200)
plt.plot([1],[1],marker='o',markersize=10, color='r')
```

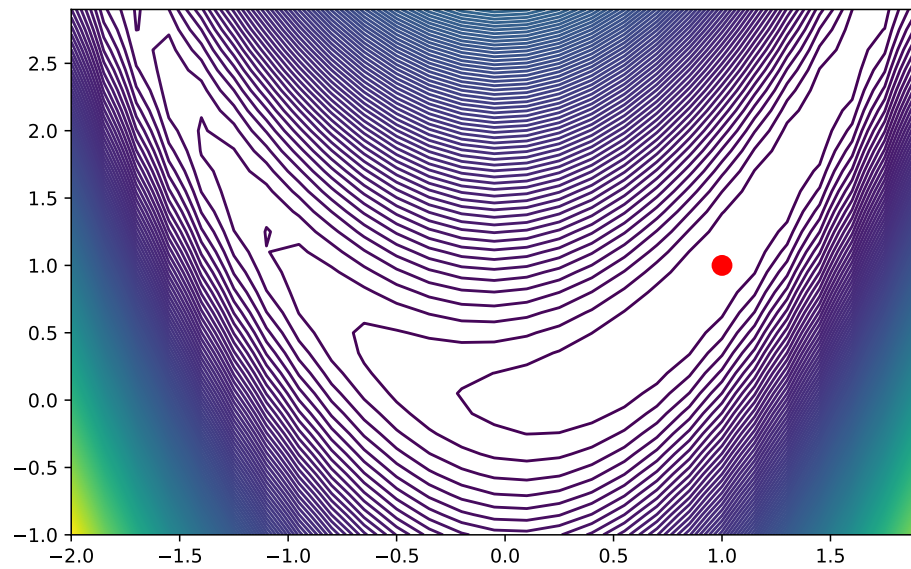


Figure 2.1.: Minumum of the Rosenbrock banana function at point (1, 1)

This minumum can be visualized in a contour chart, see figure (**banana-min?**)

In the last example we uses important information from the objective function; that is, the gradient or first derivatives. Consequently, we can use

gradient-based optimization methods such as Newton's method and conjugate gradient methods to find the minimum of this function. A potential problem arises when we do not know the gradient, or the first derivatives do not exist or are not defined. For example, imagine the following function:

$$f(x, y) = (|x| + |y|)e^{[-\sin(x^2) - \sin(y^2)]}$$

```
import math
vector = np.vectorize(np.int_)

b = 10
f = lambda x,y: (abs(x) + abs(y))*math.exp((-math.sin(x**2) - math.sin(y**2)))

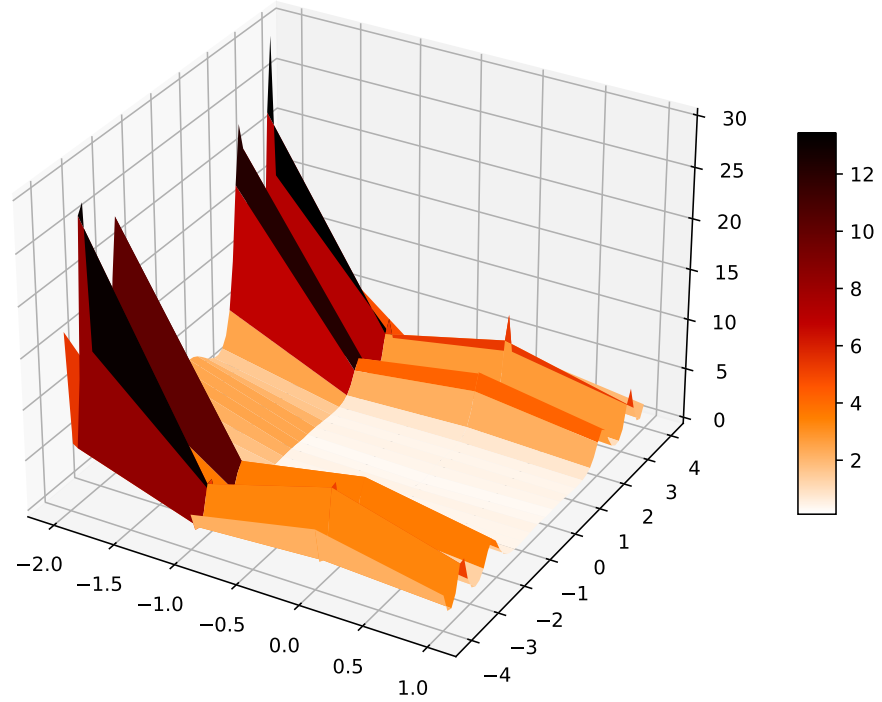
f2 = np.vectorize(f)

# Initialize figure
figRos = plt.figure(figsize=(12, 7))
axRos = figRos.add_subplot(projection='3d')

# Evaluate function
X = np.arange(-2, 2, 0.15)
Y = np.arange(-4, 4, 0.15)
X, Y = np.meshgrid(X, Y)
X = vector(X)
Z = f2(X,Y)

# Plot the surface
surf = axRos.plot_surface(X, Y, Z, cmap=cm.gist_heat_r,
                          linewidth=0, antialiased=False)
axRos.set_zlim(0, 30)
figRos.colorbar(surf, shrink=0.5, aspect=10)
plt.show()
```

2. Optimization



The global minimum occurs at $(x, y) = (0, 0)$, (**module-min?**), but the derivatives at $(0, 0)$ are not well defined due to the factor $|x| + |y|$ and there is some discontinuity in the first derivatives. In this case, it is not possible to use gradient-based optimization methods. Obviously, we can use gradient-free method such as the Nelder-Mead downhill simplex method. But as the objective function is multimodal (because of the sine function), such optimization methods are very sensitive to the starting point.

💡 Optimization Sensitive to the Starting Point

If the objective function is multimodal, and the starting point is far from the sought minimum, the algorithm will usually get stuck in a local minimum and/or simply fail.

```
plt.figure(figsize=(8, 5))
plt.contour(X,Y,Z,200)
plt.plot([0],[0],marker='o',markersize=10, color='r')
```

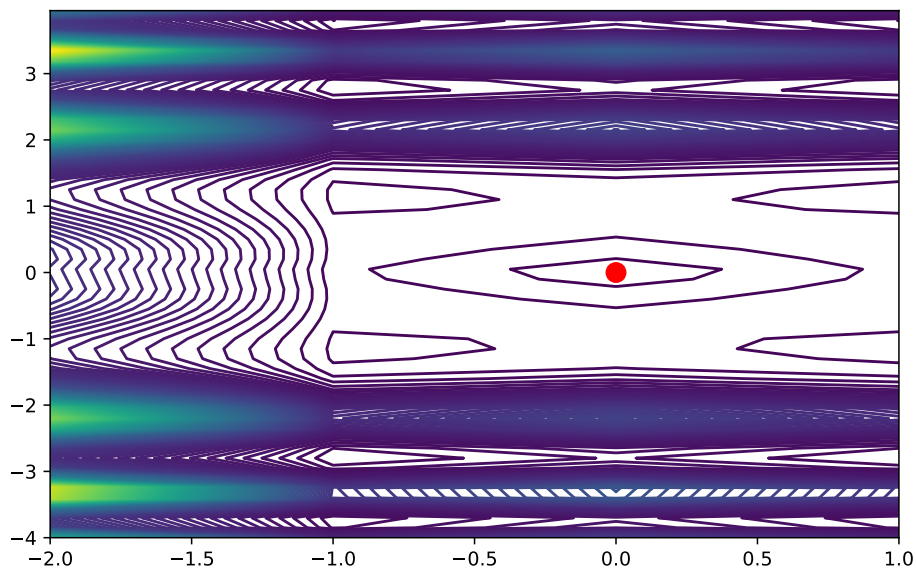


Figure 2.2.: Minumum of $f(x, y) = (|x| + |y|)e^{[-\sin(x^2) - \sin(y^2)]}$ at point (0, 0)

2. Optimization

2.1. Integer Programming

2.2. Other topics

Part III.

Supervised Learning

Supervised learning is a type of machine learning where the model is trained on labeled data, where the desired output is provided for each input. The goal of supervised learning is to learn a mapping from inputs to outputs, so that the model can make predictions on new, unseen data.

Supervised learning can be further divided into two categories: regression and classification. In regression, the output variable is continuous, and the goal is to predict a numerical value. For example, predicting the price of a house based on its square footage. In classification, the output variable is categorical, and the goal is to predict a class label. For example, classifying an email as spam or not spam.

Supervised learning algorithms can be linear or non-linear, parametric or non-parametric, and they can be based on different assumptions and mathematical models. Some examples of supervised learning algorithms are linear regression, logistic regression, decision trees, k-nearest neighbors, and neural networks.

Supervised learning is widely used in many applications, such as image and speech recognition, natural language processing, and predictive modeling. In this section, we will discuss the fundamentals of supervised learning, including the types of problems it can solve, the evaluation metrics used to measure its performance, and the algorithms and techniques used to solve those problems.

3. Linear Regression

Linear regression is a statistical method for modeling the relationship between a dependent variable (also known as the response variable) and one or more independent variables (also known as predictor variables). The goal of linear regression is to find the line of best fit (i.e., the line that best represents the relationship between the variables) through a set of data points. The equation of the line of best fit is represented by a linear equation of the form $y = mx + b$, where y is the dependent variable, x is the independent variable, m is the slope of the line, and b is the y-intercept. Linear regression can be used for both simple linear regression (one independent variable) and multiple linear regression (more than one independent variable).

$$Y_i = \beta_0 + \beta_1 \times X_i$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

to estimate the intercept β_0 we can use the following:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \times \bar{X}$$

3. Linear Regression

3.1. Learning by Examples

For this example, we will be using salary data from Kaggle. The data consists of two columns, years of experience and the corresponding salary. The data can be downloaded from [here](#).

```
/home/nmlemus/dev/ml/lib/python3.9/site-packages/IPython/core/formatters.py:
```

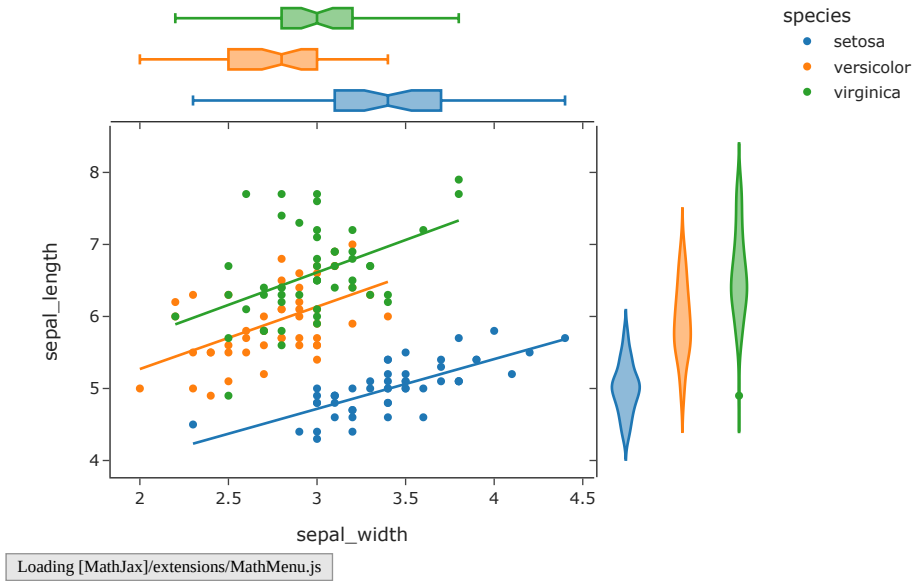
```
In future versions `DataFrame.to_latex` is expected to utilise the base impl
```

	YearsExperience	Salary
0	1.1	39343.0
1	1.3	46205.0
2	1.5	37731.0
3	2.0	43525.0
4	2.2	39891.0

```
import plotly.express as px
import plotly.io as pio
pio.renderers.default = "notebook+pdf"
pio.kaleido.scope.default_format = "png"
df = px.data.iris()
fig = px.scatter(df, x="sepal_width", y="sepal_length",
                 color="species",
                 marginal_y="violin", marginal_x="box",
                 trendline="ols", template="simple_white")
fig.show()
```

```
Unable to display output for mime type(s): text/html
```


3.1. Learning by Examples



3. Linear Regression

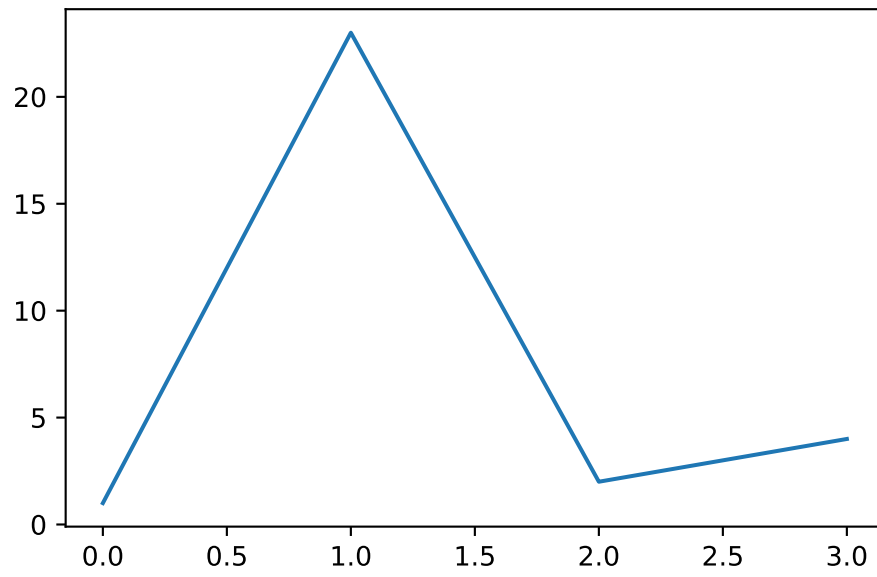


Figure 3.1.: Simple Plot

4. Logistic Regression

5. Decision Trees and Random Forests

6. Support Vector Machines

7. Neural Networks

Part IV.

Unsupervised Learning

8. Unsupervised Learning

8.1. Clustering

8.2. Dimensionality Reduction

8.3. Anomaly Detection

Part V.

Deep Learning

9. Neural Networks and Backpropagation

10. Convolutional Neural Networks (CNNs)

11. Recurrent Neural Networks (RNNs)

12. Applications in Computer Vision and Natural Language Processing

Part VI.

Advanced Topics

13. Graph Neural Networks

14. Reinforcement Learning

15. Generative Models

16. Model Interpretability

17. Conclusion

In summary, this book has no content whatsoever.

References

A. Appendices

A.1. Mathematical Proofs

A.2. Additional Resources

